

# Filtering of pulsed lidars data using spatial information and a clustering algorithm

Leonardo Alcayaga<sup>1</sup>

<sup>1</sup>DTU Wind Energy, Frederiksborgvej 399, 4000 Roskilde, Denmark

**Correspondence:** Leonardo Alcayaga (lalc@dtu.dk)

**Abstract.** Wind lidars present advantages over meteorological masts, including simultaneous multi-point observations, flexibility in measuring geometry, and reduced installation cost; but wind lidars come with the ‘cost’ of increased complexity in terms of data quality and analysis. Carrier-to-noise ratio (CNR) has been the metric most commonly-used to recover reliable observations from lidar measurements, but with severely reduced data recovery. In this work we apply a clustering technique to identify unreliable measurements from pulsed lidars scanning a horizontal plane, taking advantage of all data available from the lidars—not only CNR, but also line-of-sight wind speed ( $V_{LOS}$ ), spatial position, and  $V_{LOS}$  smoothness. The performance of this data filtering technique is evaluated in terms of data recovery and data quality, against both a median-like filter and a pure CNR-threshold filter. The results show that the clustering filter is capable of recovering more reliable data in noisy regions of the scans, increasing the data recovery up to 38% and reducing by at least two thirds the acceptance of unreliable measurements, relative to the commonly used CNR-threshold. Along with this, the need for user intervention in the setup of data filtering is reduced considerably, which is a step towards a more automated and robust filter.

## 1 Introduction

Long range scanning wind lidars are useful tools, and their adoption has grown rapidly in recent years in wind energy applications (Vasiljevic et al., 2016). Scanning wind lidars can measure time evolution and spatial characteristics of wind fields over large domains, at a lower cost of installation than meteorological masts. Nevertheless, atmospheric conditions and instrument noise can have an important impact on the data quality. For long-range scanning lidars this becomes an important issue due to the lack of additional instruments placed over the measurement area that would be useful to compare data quality, since noise can contaminate large portions of the scanning domain. The most commonly used criteria to retrieve reliable observations is a threshold on values of the Carrier-to-noise ratio, CNR, threshold that will depend on site conditions, experimental setup and the instrument manufacturer (Gryning et al., 2016; Gryning and Floors, 2019). Despite CNR threshold retrieve quality observations, its application might result in large amounts of good data rejected in regions far from the instrument, where CNR has decreased rapidly with distance. To cope with this issue Meyer Forsting and Troldborg (2016) and Vasiljević et al. (2017) have proposed filters based on the smoothness and continuity of the wind field. Such filters work by detecting discrete or anomalous steps (above a certain threshold, predefined by the user) in line-of-sight wind speed,  $V_{LOS}$ , compared to its local (moving) median. Beck and Kühn (2017) first and Karagali et al. (2018) in an adapted version, follow a different approach (here called

KDE filter, from Kernel Density Estimate) based on the statistical self-similarity of the data, which, in simple terms, means that reliable observations are alike and will be located close together in the observational space. The probability density distribution of observations (estimated via KDE) in a dynamically normalized  $V_{LOS} - \text{CNR}$  space shows that measurements likely to be valid are located in a high data density region. Observations sparsely distributed beyond a boundary defined by a threshold in  
30 the acceptance ratio, or the ratio between the probability density of any observation and the maximum probability density over the whole set of measurements, are finally identified as noise. Both approaches need the definition of one or more thresholds and a window size, either in time for the KDE filter, or in space for the wind field smoothness approach. These parameters are dependent on different characteristics of the data, like the lidar scanning pattern for instance.

Both approaches miss important and complementary information, either neglecting the strength of the signal back-scattering  
35 (quantified by CNR) or the spatial distribution and smoothness of the wind field. Moreover, in both approaches the position of observations is not taken into account, information that can shed light on areas permanently showing anomalous values of  $V_{LOS}$  or CNR, like hard targets. Including all these features within the smoothness approach is difficult, since CNR is not a smooth field like  $V_{LOS}$ . Moreover, considering smoothness and position in the KDE filter results in a computationally costly kernel density estimation, if we look for an optimal bandwidth parameter in a higher dimensional space, with a fine resolution  
40 of the kernel density estimate.

Data self similarity – over any scale in the case of fractals or a range of them in real situations (Mandelbrot, 1983) – is closely related to clustering techniques (Backer, 1995), which can classify large data sets with many different features at a relatively low computational cost. The KDE filter approach shares some characteristics with the popular *k-means* clustering algorithm MacQueen (1967), since they define one (or several for *k-means*) specific group of data belonging to an unique  
45 category (or cluster) whose size and location on the observational space will depend on data density or, more specifically, on a kernel density estimation. The main difference between these two algorithms is the way they treat sparse data points that fall in low density regions. Unlike the KDE filter, which rejects noise via the acceptance ratio, *k-means* assigns sparse points to the cluster with the nearest center, no matter if they are outliers or present unlikely values from a physical point of view.

The Density Based Spatial Clustering for Applications with Noise algorithm, or DBSCAN (Ester et al., 1996; Pedregosa  
50 et al., 2011), introduced in Section 4.3, presents several advantages over *k-means* in detecting clusters in a higher dimensional space: it introduces the notion of noise/sparsely distributed observations, it does not need prior knowledge of the number of clusters in the data and it is capable of identifying clusters of arbitrary shape. To the best of our knowledge, this is the first time that this type of clustering algorithm is applied to identify not reliable observations from pulsed lidars. This approach, which can be understood as a natural extension of the KDE filter, is compared to the smoothness based filter on two types of data:  
55 synthetic wind fields data as a controlled test case, and real data.

This paper is organized as follows: Section 2 describes the real data used to test the different filtering approaches, and Section 3 presents the synthetic data used during a controlled test as well as the methodology to obtain it. Section 4 then gives a description of the different filters applied in this study to both data sets, to continue with the definition of the performance tests in Section 5. In Section 6 the performance tests are presented along with a discussion on their validity and quality. Section 7

**Table 1.** Characteristics of the Balconies experiment, from Karagali et al. (2018). The scans are not instantaneous neither totally synchronous, with a horizontal sweep speed of  $2^\circ/\text{s}$  in the azimuth direction in a range of  $90^\circ$ , with a total time of 45 s per scan.

Phase	Measurement start	Measurement end
50 m a.g.l. (1)	2016-04-12 12:45:41	2016-06-17 12:48:01
200 m a.g.l. (2)	2016-06-29 13:35:56	2016-08-12 09:09:55
Scanner	Location coordinates, [m]	Scanning pattern, west
Southern lidar	492768.8 (East) 6322832.3 (North)	$344^\circ$ - $256^\circ$ , $2^\circ$ steps
Northern lidar	492768.7 (East) 6327082.4 (North)	$196^\circ$ - $284^\circ$ , $2^\circ$ steps

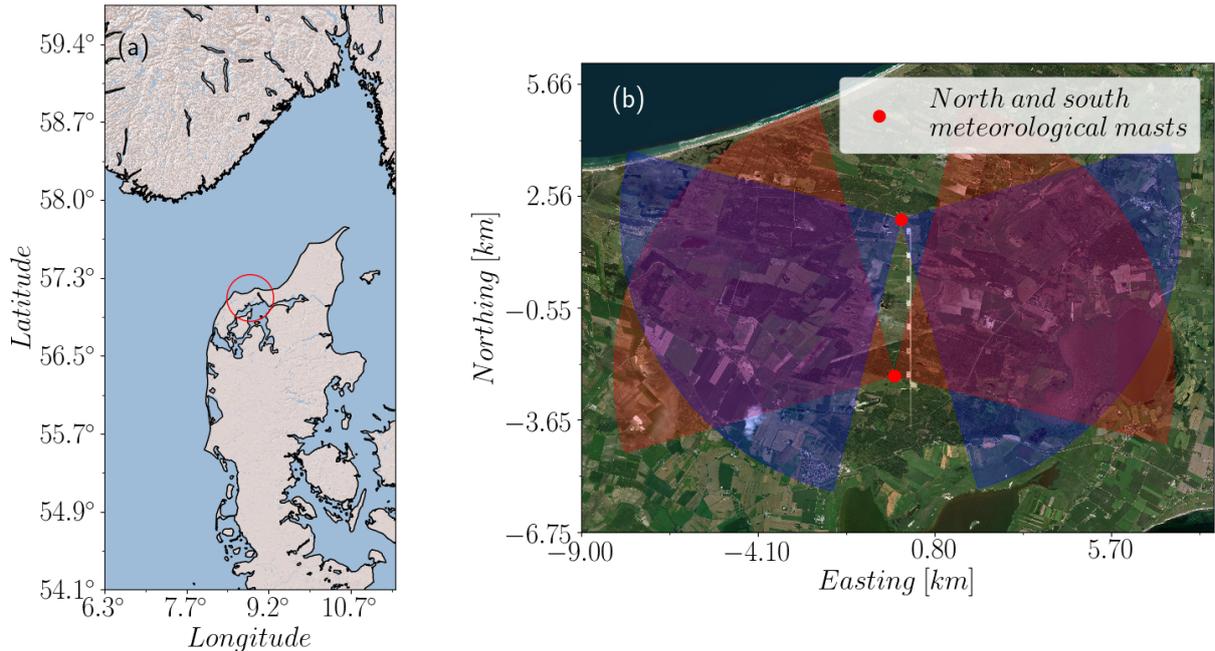
60 discusses the quality of the methodology behind the tests and the advantages and disadvantages of the proposed approach. Section 8 presents the conclusions of this study.

## 2 Real data: Østerild Balconies experiment

The filtering techniques presented here were tested on lidar measurements made at the Østerild Test Centre located in northern Jutland, Denmark, see Figure 1. Known as the Østerild Balconies experiment (Mann et al., 2017; Karagali et al., 2018; Simon and Vasiljevic, 2018), this measuring campaign aimed to characterize horizontal flow patterns above a flat, heterogeneous forested landscape at two heights relevant for wind energy applications, covering an area of around  $50 \text{ km}^2$ , and a wide range of scales, both in time and space.

The experiment consist of two measuring phases (see Table 1) with two long-range WindScanners performing Plan Position Indicator (PPI) scanning patterns, aligned in the North-South axis and installed at 50 m a.g.l. during phase 1 and 200 m a.g.l. in phase 2. WindScanners (Vasiljevic et al., 2016) consist of two or more spatially separated lidars which are synchronized to perform coherent scanning patterns, allowing the retrieval of two or three dimensional velocity vectors at different points in space. These experiments were conducted between April and August of 2016 (Simon and Vasiljevic, 2018). In each phase, the northern and southern lidars scanned in the West and East direction relative to the corresponding meteorological masts, where they were installed. The data used in this study originated from both phases of the experiment, with PPIs pointing to the west. For more details about the experiment, lidars and terrain characteristics see (Karagali et al., 2018; Vasiljevic et al., 2016; Simon and Vasiljevic, 2018).

This dataset is well suited to test different data filtering techniques. A large measurement area will be affected by local terrain and atmospheric conditions, like clouds or large hard targets. Moreover, at this scale lidars reach their measuring limitations, since the back-scattering from aerosols decrease rapidly with distance (Cariou, 2015).



**Figure 1.** (a) Location of the Østerild Turbine Test Center, place of the Balconies experiments, northern Jutland, Denmark (copyright 2009 Esri). (b) Detail of the test center site, with the location of the meteorological masts where north (blue) and south (red) WindScanners were installed. During the measurement campaign the PPI scans pointed both west in some periods and both east in other (copyright 2017 DigitalGlobe, Inc.).

### 80 3 Synthetic data

Assessing and comparing the performance of filters is challenging with no reference available to verify that rejected or accepted observations are reliable or bad observations. This is especially difficult for long-range scanning lidars, since their measurements cover large areas and, due to spatial variability, a valid reference would need several secondary anemometers scattered over the scanning area. Testing filters on a controlled and synthetic data set, contaminated with a well defined noise, presents an option to deal with this problem. In this study, the filters presented in Section 4 are tested on individual scans sampled from synthetic wind fields generated using the Mann turbulence spectral tensor model (Mann, 1994), and contaminated with procedural noise (Perlin, 2001).

#### 3.1 Synthetic wind fields generation

Synthetic PPI scans are sampled by a lidar simulator from synthetic wind fields generated via the Mann-model (Mann, 1998) in a horizontal, two-dimensional square domain of 2048 x 2048 grid points, with dimensions 9200 m x 7000 m. The generated turbulence fields are the result of input parameters of the of turbulence spectral tensor model, namely, length-scale,  $L$ , turbulence

**Table 2.** Synthetic wind field characteristics and parameters.

Parameter	Values
$L$ , m	62, 125, 250, 500, 750, 1000
$\alpha\epsilon^{2/3}$ , $\text{m}^{4/3}\text{s}^{-2}$	0.025, 0.05, 0.075
$\Gamma$	0, 1, 2, 2.5, 3.5
Number of seeds used	10
Mean wind speed, $U$ m/s	15
Mean wind speed direction range, degrees	90 to 270
Total number of scans generated	4305

energy dissipation  $\alpha\epsilon^{2/3}$ , and anisotropy,  $\Gamma$ . The fields generated correspond to wind speed fluctuations, to which the desired average wind speed mean is subsequently added. Depending on the initial random seed used, different wind field realizations with the exact same turbulence statistics can be generated. For details on wind field generation using the the Mann-model, refer to Mann (1998). Table 2 shows the range of values used for the generation of two-dimensional wind fields. Large values of  $\alpha\epsilon^{2/3}$  or small scale turbulence for instance, mean that sudden spatial changes in wind speed are more likely, which increase the false identification of outliers. Mean wind direction, turbulence anisotropy and length scale will also affect the sampling due to the lidars measuring characteristics.

### 3.2 Lidar simulator

Lidar simulators has been presented previously by Stawiarski et al. (2013) and Meyer Forsting and Troldborg (2016). They sample  $V_{LOS}$  values from wind fields generated via Large Eddy Simulations (LES), mimicking the operational principle of lidars by proper time and spatial (probe volume) averaging of the background wind field. The lidar simulator presented here follows the same principles, this time sampling from synthetic wind fields generated via Mann-model.

The simulator receives scanning pattern characteristics as input (beam range, range gate step, azimuth angles range and azimuth angle steps) to generate a primary mesh with the sampling positions on top of background wind field. Following the measuring principle of the lidar, the  $V_{LOS}$  observed at each position in this mesh will represent averages of a continuous along each range gate step (due to probe volume averaging) and an average of many azimuth positions within the azimuth step, due to the almost continuous sweep of the lidar's beam. The simulator mimics this generating a secondary, refined mesh with  $N_r$  points in each range gate and  $N_\phi$  beams within each azimuth step. The background wind field components,  $U$  and  $V$ , are then interpolated on this secondary mesh and projected on each refined beam to obtain  $V_{LOS}$  using equation (1), with  $\theta$  being the corresponding beam azimuth angle.

$$V_{LOS} = \cos(\theta)U + \sin(\theta)V \quad (1)$$

**Table 3.** The characteristics of the lidar simulator and real long-range lidar (Karagali et al., 2018; Vasiljevic et al., 2016) used for the controlled test of the filters.

	<b>Simulator</b>	<b>Real</b>
Azimuth range	256° - 344°	256° - 344°
Azimuth step	2°	2°
Beam length	7000 [m]	7000 [m]
Range gate length, $\Delta p$	35 [m]	35 [m]
Full width at half maximum, $\Delta l$	75 [m]	75 [m]
Sweeping time per scan	Instantaneous	45 [s]
Primary mesh size (radial x azimuth)	45 x 198	-
Secondary mesh size at each range gate ( $N_r \times N_\phi$ )	21 x 51	-
Total secondary mesh size ( $N_r \times N_\phi$ )	21 x 51	-

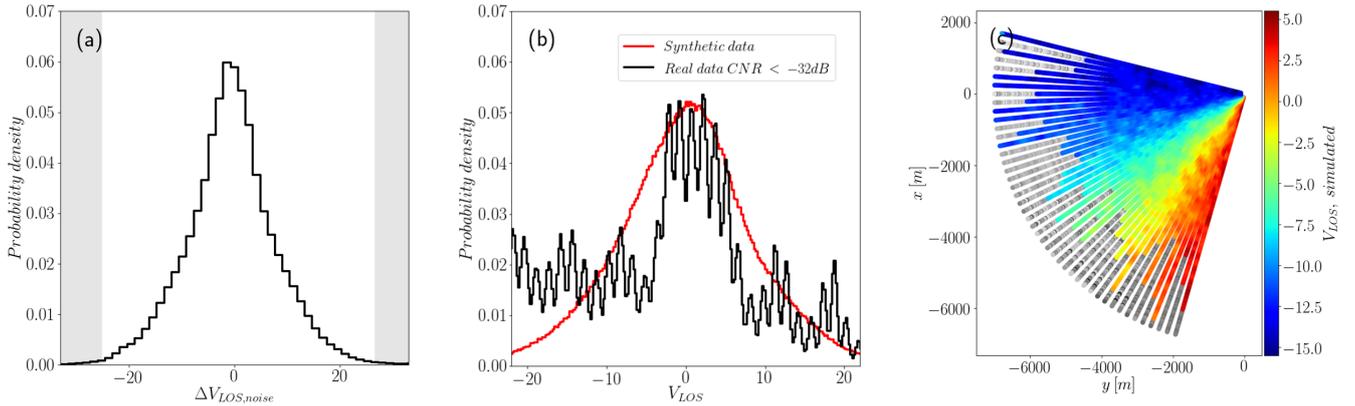
The final step is the spatial (probe volume) averaging, and the azimuth (sweeping) averaging around each position in the primary mesh. Spatial averaging is done applying a weighting function on all  $V_{LOS}$  along each refined beam. The weighting function used here is defined in equation (2), as in Banakh and Smalikhov (1997) and Smalikhov and Banakh (2013). This function will assign weights to each point in the refined beam according to its distance to the range gate position in the primary mesh,  $F$ , and the instrument probe volume parameters, namely, range gate length,  $\Delta p$ , and full width at half maximum,  $\Delta l$  (cf. Table 3). Here,  $\text{Erf}(x)$  is the error function, and  $r_p$  the beam width contribution to the volume averaging.

$$w = \frac{1}{2\Delta p} \left\{ \text{Erf} \left[ \frac{(r - F) + \Delta p/2}{r_p} \right] - \text{Erf} \left[ \frac{(r - F) - \Delta p/2}{r_p} \right] \right\}; \quad r_p = \frac{\Delta l}{2\sqrt{\ln(2)}} \quad (2)$$

The azimuth averaging is the arithmetic mean of the  $N_\phi$  values of  $V_{LOS}$  at each range gate after spatial averaging. It represents the accumulation information of the back-scattered signal spectra as they sweep an azimuth sector before estimation of the spectral peak and  $V_{LOS}$ .

### 3.3 Synthetic noise generation

The most simple noise that can be used to contaminate synthetic scans is sparse, uniformly distributed outliers. This noise, also known as salt and pepper noise, is easily detected and eliminated by median-like filters, when extreme discrete steps affect the smoothness of an image (Huang et al., 1979; Burger and Burge, 2008). Nevertheless, noise in real scans comes as regions of anomalously high and/or low  $V_{LOS}$  and they can pass through the filter undetected. Procedural noise, introduced by (Perlin, 2001) to recreate synthetic textures on surfaces for computer graphics applications, creates regions of coherent noise that resembles better the spatial distribution of scanning lidars measurements. For the two-dimensional case, the procedural noise function  $N(x, y)$  maps two-dimensional coordinates,  $(x, y)$ , onto the range  $[-1, 1]$  as follows,



**Figure 2.** Procedural noise on synthetic scans. (a) Distribution of  $\Delta V_{LOS,noise}$ , the noise added. Maximum values are within the observable range between  $[-35, 35]$  [m/s]. (b) Distribution of real  $V_{LOS}$  with low CNR values (black) and contaminated, synthetic  $V_{LOS} + \Delta V_{LOS,noise}$  (red) for a mean wind direction facing the scan. (c) Individual scan showing the increasing fraction of added noise (grey) with distance.

135

- A two-dimensional grid of  $m$  by  $n$  elements is generated, and a pseudo-random, two-dimensional unit gradient,  $\mathbf{g}_{ij} = (g_x, g_y)$ , is assigned to each grid point  $(x_i, y_j)$ . The pseudo-randomness rises from the fact that  $\mathbf{g}_{ij}$  are picked from a pre-computed list of gradients with length  $l \ll m \times n$ . We select values from this list using the index permutation grid  $p_{ij} \in \{0, \dots, l\}$  also with  $m \times n$  elements. Then,  $\mathbf{g}_{ij}$  will correspond to the gradient in the position  $p_{ij}$  of the pre-computed list. Elements  $p_{ij}$  are shuffled for each realization.
- For each grid point  $(x_i, y_j)$  enclosing  $(x, y)$ , a distance vector  $\mathbf{d}_{i,j} = (x - x_i, y - y_i)$  is generated.
- Finally, the noise function is the sum of dot products,  $N(x, y) = \sum_q w_q (\mathbf{g}_{ij}^q \cdot \mathbf{d}_{i,j}^q)$ , for  $q$  grid points surrounding  $(x, y)$ . Weights  $w_q$  correspond to  $w_q = C \frac{1}{\|\mathbf{d}_{i,j}^q\|}$ , and  $C$  a normalization constant to ensure that  $N(x, y) \in [-1, 1]$ .

140

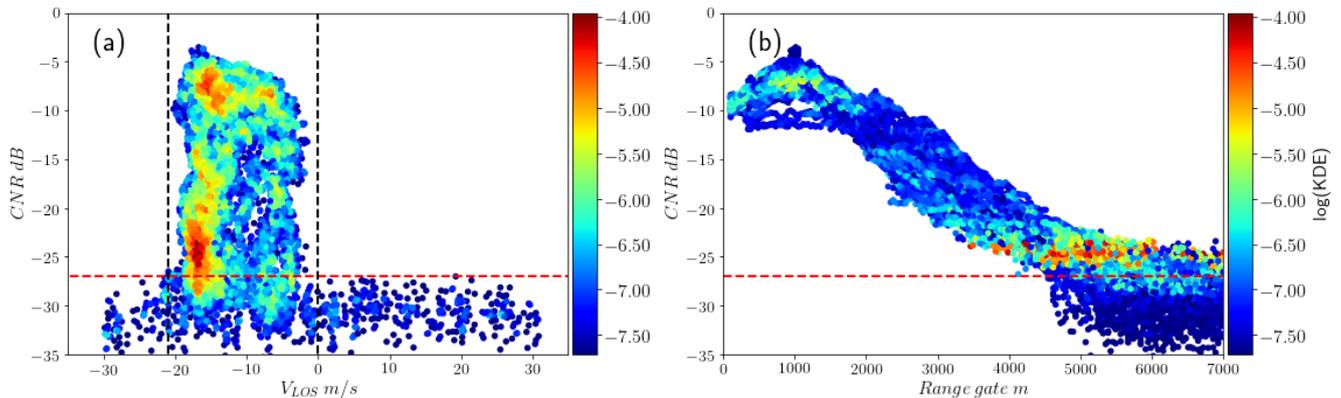
The function  $N(x, y)$  allows the generation of noisy regions, than can be distributed according to back-scatter decay with distance. Three bands centered at 50% , 70% and 90% of the total beam length (and spanning over the entire azimuth range) have an increasing fraction of noise, contaminating the 30%, 60% and 90% of the observation, respectively. The noise amplitude is 35 [m/s], the limit of the observable range for the instruments described in Section 2. Figure 2 (c) show one contaminated scan and its increasing contaminated area as we move along the beams. The same Figure shows the distribution of the noise generated by the algorithm after scaling, and the probability distribution of contaminated synthetic  $V_{LOS}$  compared to real data

145

with low values of CNR. The distribution of real data presents heavier tails than the ones generated, with higher probability of observing extreme values of  $V_{LOS}$ . Modeling real noise is difficult, since the process that generates it depends on the measuring principle of the lidar and atmospheric conditions. The synthetic noise used here does not intend to be totally realistic, but more subtle and smoother than the one observed in real measurements, making the identification of contaminated points more difficult.

#### 4.1 CNR threshold

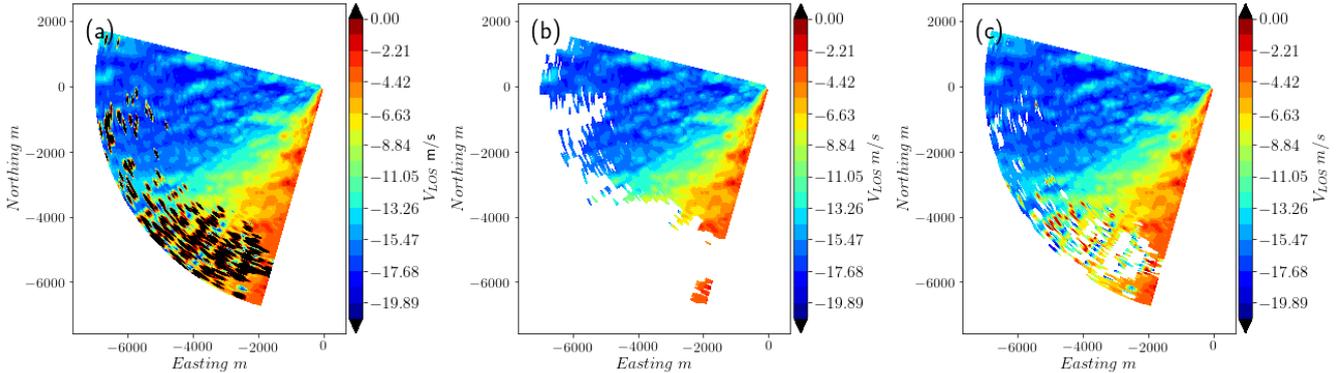
CNR thresholds are well known and lidar manufacturers usually recommend values for rejection of signals with poor backscattering or hitting hard targets (Cariou, 2015). However, the selection of an appropriate threshold for CNR that assures data quality and good data recovery is not easy. Figures 3 and 4 show data from a scan with noisy observations from CNR values below  $-27$  dB. Both, extreme and limited values of  $V_{LOS}$ , show low CNR values in the distant region of the scan, and data loss results after the application of the CNR threshold (Figure 4 (b)). When a limit to  $V_{LOS}$  is applied instead, Figure 4 (c) show that the smoothness in  $V_{LOS}$  is lost in the lower part of the scan. A conservative threshold of  $-24$  dB is used here, since the resulting  $V_{LOS}$  probability distribution show very little outliers and it can be used as a reference when the performance of the filters proposed are compared.



**Figure 3.** (a) CNR and  $V_{LOS}$  for one scan from Balconies experiment, including the probability density (KDE). Observations with  $CNR > -27$  dB (dashed red line) show a limited range of  $V_{LOS}$  (dashed black line). A portion of observations with high probability density remain in the rejection area. (b) CNR v/s distance for the same data. Observations with low CNR values and high probability density can be found in the distant region of the scan.

#### 160 4.2 Median-like filter

The median filter arise as a viable option for detecting erroneous measurements, since it is well known that this type of non-linear filter is suited to detect and filter noise that present distributions with large tails. Here we use an adaptation of the traditional median filter used in the image-processing community, closely related to the three-stage filtering technique described in Menke et al. (2019): observations are not replaced by the local moving median but excluded if the absolute difference between their value and the local moving median is above a certain threshold,  $\Delta V_{LOS, threshold}$ . Unlike Huang et al. (1979), The two-dimensional moving window is replaced by a two one-dimensional window instances, the first in line-of-sight



**Figure 4.** (a) Un-filtered scan with  $V_{LOS}$  values outside the range  $[-21, 0]$  in Figure 3 (a) in black. (b) Filtered scan with  $CNR > -27$  and the resulting data loss in the upper part of the scan. (c)  $V_{LOS}$  within the  $[-21, 0]$  range, showing anomalous values in the lower part of the scan.

or radial direction,  $r$ , and finally in the azimuth direction,  $\theta$ , considering the polar coordinates of the scan. This simplification reduces the computation time importantly.

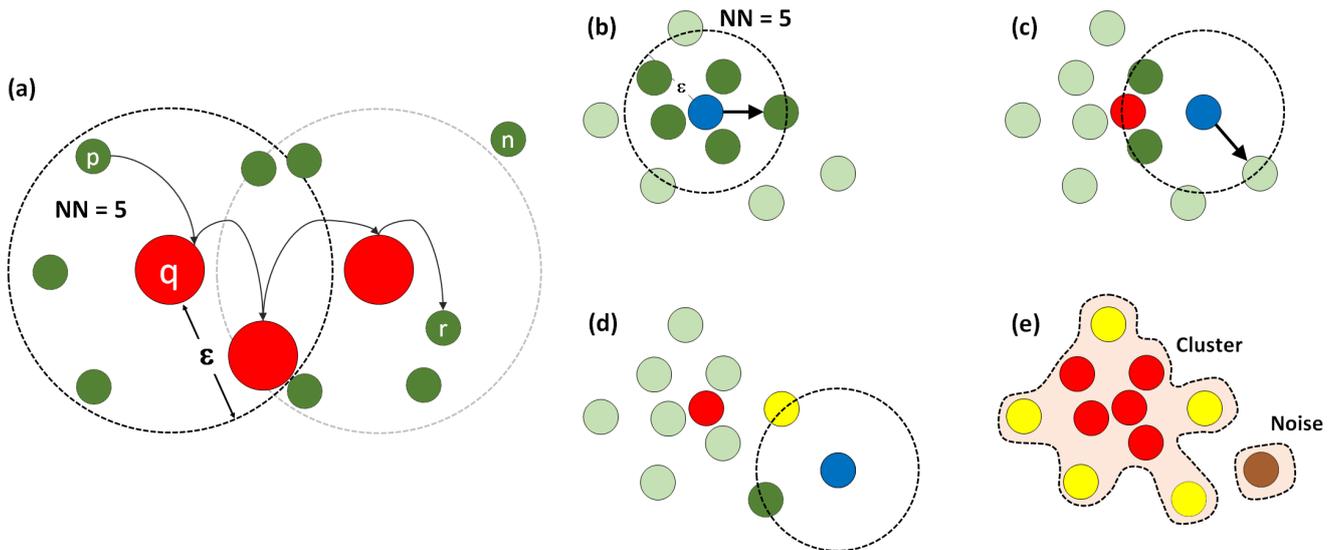
The input parameters of this filter will be the size (or number of elements) of the moving windows in the radial and azimuth directions,  $n_r$  and  $n_\phi$  respectively, and  $\Delta V_{LOS, threshold}$ . For fixed values of  $\Delta V_{LOS, threshold}$ ,  $n_r$  and  $n_\phi$ , the spatial structure of wind speed fluctuations will have an important effect on the recovery rate and noise detection of this filter. A sensitivity analysis carried out using the metrics presented in section 5.1 on the synthetic data set, shows that the optimal set of parameters is  $n_r = 5$ ,  $n_\phi = 3$  and  $\Delta V_{LOS, threshold} = 2.33$  m/s (See Appendix A). This set is used both for artificial and real data. The filter does not include a time window, and it is applied on individual scans.

### 175 4.3 Filtering using a clustering algorithm

If we represent lidar observations as  $m$ -dimensional vectors, with  $m$  the number of features/parameters of the data, measurements not affected by poor back-scattering or noise will cluster together in regions of high data density, as shown in Figures 4 and 3. The approach presented here identifies such clusters applying DBSCAN on data described by  $CNR$ ,  $V_{LOS}$  and, additionally, spatial location and smoothness features, which help to make clusters more distinguishable.

180 DBSCAN identifies clusters and noise based on two parameters: a neighbourhood size,  $\epsilon$ , and a minimum number of nearest neighbours,  $NN$ . The parameter  $\epsilon$  is the euclidean distance from one observation to the limits of a neighborhood that might contain  $NN$  (or more) nearest neighbors. Intuitively, these parameters will define the minimum density that a group of data points needs to have to be identified as a cluster. Observations within a cluster fall into the following categories,

- Core point: points  $q$  whose  $\epsilon$ -neighborhood contains  $NN$  or more points.
- 185 – Direct density reachable point: points  $p$  which are reachable by  $q$  by laying within its  $\epsilon$ -neighborhood.

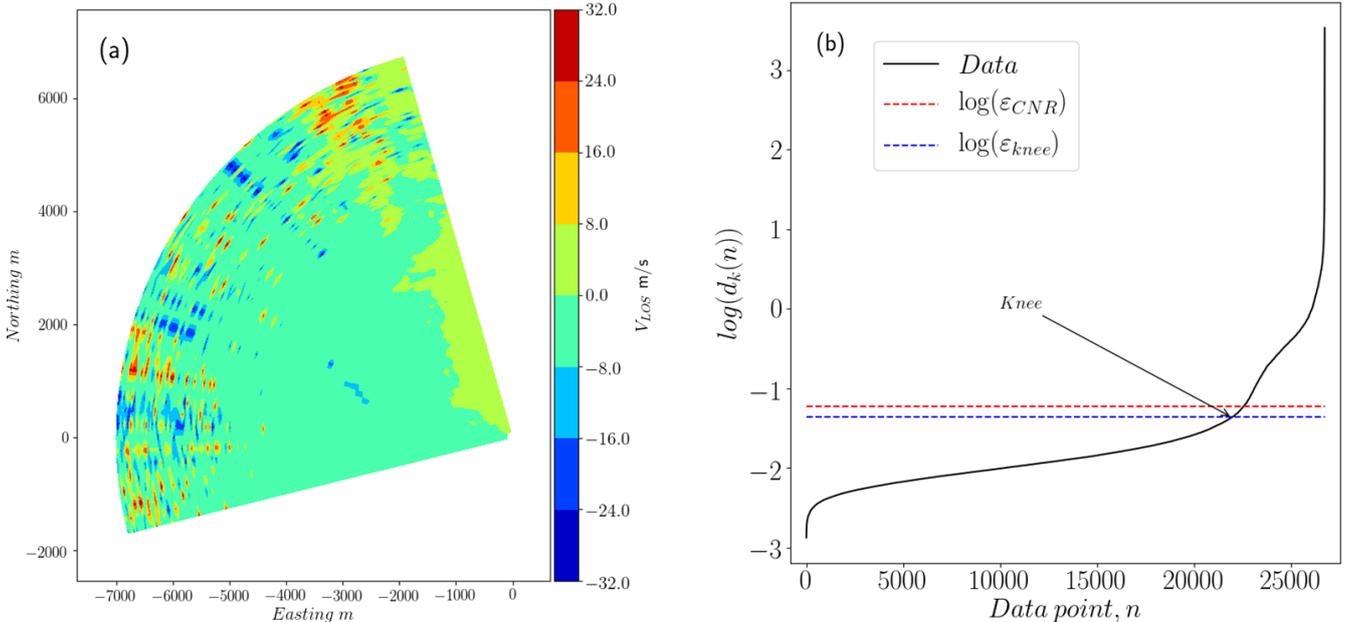


**Figure 5.** (a) DBSCAN algorithm definitions: direct density reachable point  $p$  (reachable by the core point  $q$ ) and density reachable and density connected points  $p$  and  $r$ . Here point  $n$  does not belong to any of these categories, but noise. The DBSCAN algorithm working: (b) The current point being evaluated has the minimum number of nearest neighbours required,  $NN$ , within a neighborhood of size  $\epsilon$ , classified as a *core point* (red) (c) The next point has less than  $NN$  neighbours, but one of them is a core point and becomes a *border point* (yellow) (d) A point with neither  $NN$  neighbours, nor core points within  $\epsilon$ , classified as *noise* (brown) (e) The final cluster and noise. The former is a collection of density connected points.

- Density reachable point: points  $p$  reachable by a point  $r$  through one or a set of directly connected core points  $q$ .

DBSCAN travels across data points identifying core points, border points (density reachable points with at least one core point within the  $\epsilon$ -neighbourhood) and noise, or points that do not belong to any of the categories described above. Finally, the algorithm separates clusters as individual groups of density connected points. Figure 5 shows schematically these definitions and how the algorithm works.

The input parameters  $\epsilon$  and  $NN$  have a significant influence on the number and characteristics of the clusters detected. For example, large  $\epsilon$  together with a small  $NN$  will end up with sparse clusters that might include noise. In order to find the parameters separating the least dense cluster from noise, we fix  $NN$  to a certain value  $k$  and determine  $\epsilon$  from the data density distribution. The latter is well described by the  $k$ -distance function,  $d_k(n)$ , which represents the distances from all data point  $n$  to their respective  $k$ -th nearest neighbour, sorted in ascending order. When  $k$  is 5 for instance,  $d_5(n)$  in Figure 6 shows sudden changes (or knees) that give some indications about the data density distribution. The knee highlighted represents a limit between a group of reliable observations and the one growing fast towards noisy data. The positions of these knees in the graph correspond to the peaks in the curvature of the  $d_k(n)$ ,  $\kappa(n)$  in expression (3). In this expression primes correspond to the derivatives of  $d_k(n)$  with respect to a  $n$ . The continuous version of  $d_k(n)$  is made by spline-fitting on a reduced set of



**Figure 6.** (a) Scan from the Balconies experiment (phase 1) with a 48% of data points in the range of reliable observations with  $CNR \in [-24, -8]$  dB (b) Logarithm of sorted distances to the 5-th nearest neighbour for each point in a data set. The total number of observations corresponds to three consecutive scans, or 26730 points. The sorted 5-th distances show three knees separating three types of structures: reliable observations with distances below  $\varepsilon_{knee}$ , an overlapping region where the distance between points grows faster, and pure noise or non structured data.

200 uniformly distributed points over the original data set.

$$\kappa = (d_k(n))'' / (1 + (d_k(n))')^2)^{3/2} \quad (3)$$

When scans are very noisy, the selection of a proper value of  $\varepsilon$  is difficult, since knees are located closer together and a larger fraction of observations show a fast growing  $d_k(n)$ , as expected. In this case, the fraction of points showing a reliable CNR values is taken into account and  $\varepsilon$  is estimated by expression (4). Here  $f_{CNR}$  corresponds to the fraction of observations  
 205 CNR values within the range  $[-24, -8]$  and the constants  $c_1$  and  $c_2$  are obtained from the upper and lower bounds of  $\varepsilon$  in the data, respectively.

$$\varepsilon_{CNR} = c_1 f_{CNR} + c_2 \quad (4)$$

The set of features considered when filtering synthetic data does not include CNR, because it is not available from the the lidar simulator described in Section 3. For synthetic and real data sets we consider spatial location (azimuth and radial  
 210 positions) and smoothness as additional features. The latter,  $\Delta V_{LOS}$ , corresponds to the median difference in  $V_{LOS}$  between a specific position and its direct neighbours in one individual scan.

Since we consider features that vary importantly in magnitude (CNR and range gate distance for instance), we normalize the data before the application DBSCAN. This step is necessary for the estimation of meaningful distances between observations, basis of this approach. There are several ways to do this. Here, the data in each feature is centered by subtracting its median, and scaled according to its inter-quantile range. This aims to minimize the influence from outliers in the normalization.

The clustering filter is implemented to be a non-supervised classifier, and does not need more input parameters than the different features and the number of scans put together as a batch before filtering. The latter is set to three in this case, to speed up calculations and avoid creating clusters from noisy regions. From this point of view, this filter is also dynamic as that of Beck and Kühn (2017) when applied to a real data set, since it will consider the data structure within a period limited to 135 seconds (3 scans of 45 seconds in our case), and characteristics of temporal evolution of the data is indirectly taken into account. For the synthetic data used in this test, more than one scan filtered per iteration gives enough data density in noisy and reliable areas of the observational space. We speculate that scans that are correlated in time will enhance the self-similarity of the data, thus improving the performance of the filter. Turbulence structures with length scales in a range between the range gate size and the scanning area size will evolve at a slower rate than the time elapsed between consecutive scans.

## 5 Performance metrics

### 5.1 Synthetic data

Expressions (5) to (7) defines three metrics to assess the performance of the filters, given prior knowledge of the position and magnitude of noise in a controlled case with  $N$  observations. The fraction of noise detected,  $\eta_{noise}$ , quantify the relative importance of true positives, or the difference between observations identified as noise,  $N_{noise}$ , and false positives,  $N_{pos}$ , over the total number of contaminated observations. The fraction of good observations recovered,  $\eta_{recov}$ , give an idea of the true negatives over the total number of non-contaminated observations,  $N_{non-cont}$ . True negatives are not not equal to  $N - N_{noise}$ , since the latter might include false negatives,  $N_{neg}$ . The relative importance of this two metrics, for a given fraction of noise in a contaminated scan,  $f_{noise}$ , is quantified by  $\eta_{tot}$ , which takes into account cases with a large fraction of noise detected and low recovery rate, and vice-versa.

$$\eta_{noise} = \frac{N_{noise} - N_{pos}}{N_{cont}} \quad (5)$$

$$\eta_{recov} = \frac{N - (N_{noise} + N_{neg})}{N_{non-cont}} \quad (6)$$

$$\eta_{tot} = f_{noise}\eta_{noise} + (1 - f_{noise})\eta_{recov} \quad (7)$$

In the absence of reference measurements, the quality of the data retrieved after filtering is assessed by comparing the distribution of radial wind speeds for very reliable observations (with CNR values within the range between -24 to -8 dB) with the distribution of filtered observations that fall out of this range. Observations out of the reliable range population usually show a probability density function (or PDF) with heavier tails, like the pdfs in Figure 7. Here we understand a heavy tailed pdf as a distribution that slowly goes to zero and show higher probability density for values beyond the  $3\text{-}\sigma$  limit (or 3 standard deviation limit), when compared to the normal distribution, evidence of a higher probability of occurrence of outliers or extreme values. The recovering rate of observations beyond the [0.003, 0.997] quantile range of the reliable  $V_{LOS}$  (shaded area in Figure 7) could shed information about the quality of the data retrieved by the filter.

Other metric is the similarity between pdf of reliable and non reliable data, after filtering. The distance between both probability density functions can be compared with similarity metrics like the Kolmogorov-Smirnov test (Kolmogorov, 1933) or Kullback–Leibler (KL) divergence (Kullback and Leibler, 1951). The former test measures the statistical similarity between two random variables,  $X_1$  and  $X_2$ , by estimating the statistical distance,  $D$  (or K-S statistic), between their cumulative distribution functions,  $F_1(x)$  and  $F_2(x)$ , as the supreme of their difference,

$$D_K = \sup_x \|F_1(x) - F_2(x)\| \quad (8)$$

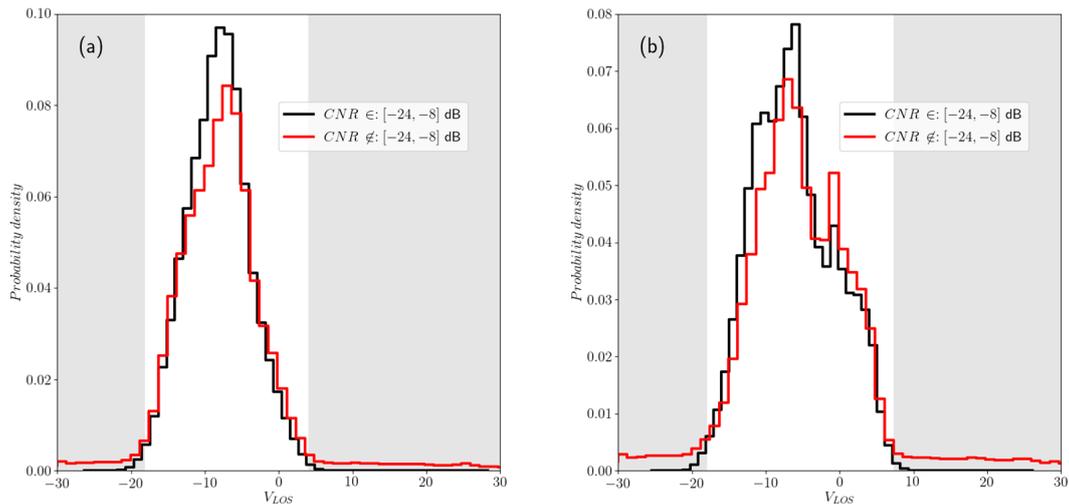
The null hypothesis here is that two realizations are from the same distribution, if the K-S statistic is such that its two tailed  $p$ -value is above a certain level  $\alpha$ . Due to the amount of data analyzed here is huge—we analyzed over 20000 scans for the two phases of the Østerild campaign, each with 8910 data points, over almost 10 days—this similarity test is very precise, but also very strict rejecting the null hypothesis for small deviations between  $F_1(x)$  and  $F_2(x)$ . Nevertheless, the K-S statistics can be used to compare which probability distribution after filtering is closer to the one representing the reliable observations.

The KL divergence is a measure of similarity, or overlapping of two distributions  $P_1$  and  $P_2$ , with realizations  $X_1$  and  $X_2$ , respectively. It is used in different applications to shed light on the loss of information when  $X_1$  is represented by  $P_2$  or vice-versa and is defined by the expression (9).

$$D_{KL} = \sum_x P_1(x) \log \left( \frac{P_2(x)}{P_1(x)} \right) \quad (9)$$

Both metrics will be used to estimate how the distribution of non reliable observations of  $V_{LOS}$  is modified after filtering, and if the new distribution is similar (or close, in a statistical distance sense) to the probability density of reliable observations of the radial wind speed, shown in Figure 7 for phases 1 and 2 of the measurement campaign, respectively.

Both performance metrics, the recovery rate of abnormal measurements in the tails of the pdf of reliable observations and its statistical distance to the pdf of filtered non reliable observations, will be assessed for the median-like filter, the clustering filter and also for data filtered with a CNR threshold of -29 dB, following (Gryning and Floors, 2019).



**Figure 7.** Probability density function of reliable observations of  $V_{LOS}$  (black solid line) and non reliable observations (red, solid line) for (a) Phase 1 of Balconies experiment with scans performed at 50 m a.g.l. and (b) Phase 2 of the same campaign, with scans performed at 200 m a.g.l.

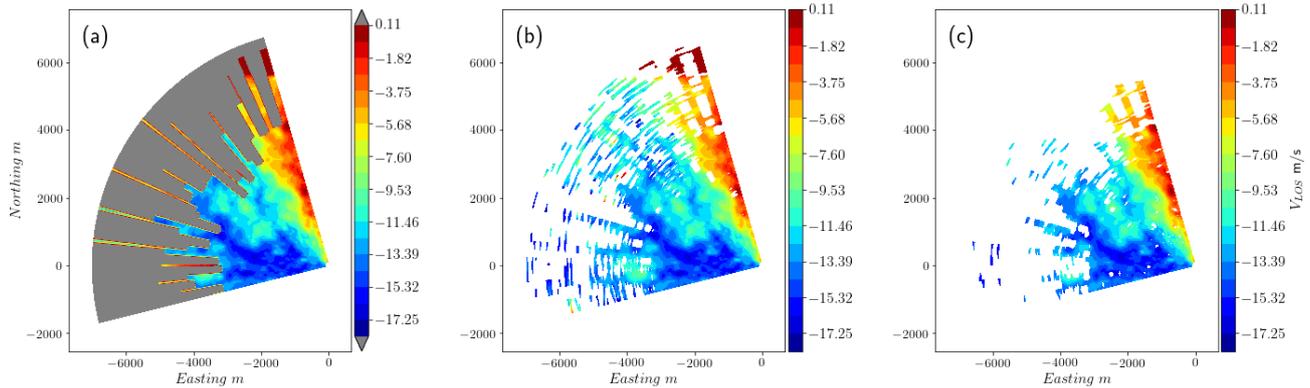
## 270 6 Results

### 6.1 Synthetic data

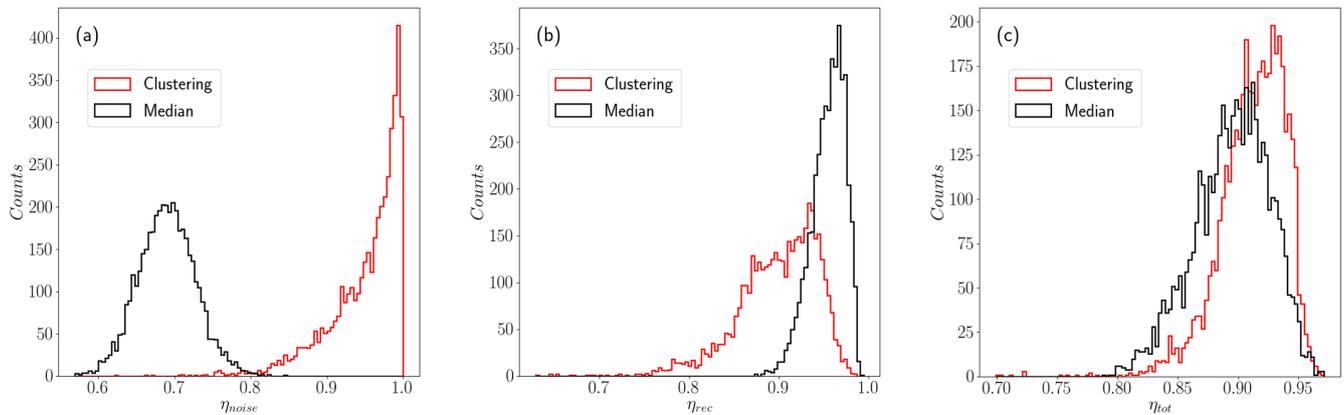
In Figure 8 we can see the result of the two filters applied on one synthetic scan contaminated with procedural noise. The contaminated observations are indicated by the grey area in this scan. Extreme values contaminating  $V_{LOS}$  are identified by both filters without problems, but subtle alterations on the original values of the scan are hard to detect for the median-like filter.

275 The clustering filter performs very efficiently detecting this type of contaminated observations and filters almost all the noise. Both filters repeat this behavior in all the synthetic scans used for this controlled test, as can be seen in Figure 9, which shows the resulting metrics of the two filters applied on the whole synthetic data set. Looking at  $\eta_{tot}$ , both filters show similar mean values and spread, with the clustering filter performing slightly better. The difference becomes noticeable when we see  $\eta_{noise}$ , which for the clustering filter show a mean value of 0.95, far larger than the 0.67 of the median-like filter. The latter result  
280 could be problematic if the median-like filter is used, since noise contaminating the filtered scan will result in non realistic wind fields after reconstruction.

Both filters perform well when evaluated in terms of  $\eta_{rec}$ , with the median-like filter showing a higher mean fraction of good observations retrieved, 0.96, compared with the 0.89 of the clustering filter. This result is expected, since the median-like filter is more permissive regarding fluctuations that can seem locally anomalous for the clustering filter.



**Figure 8.** (a) Contaminated synthetic scan with noise indicated by grey area. (b) Scan filtered using the median-like approach. (c) Clustering filter.

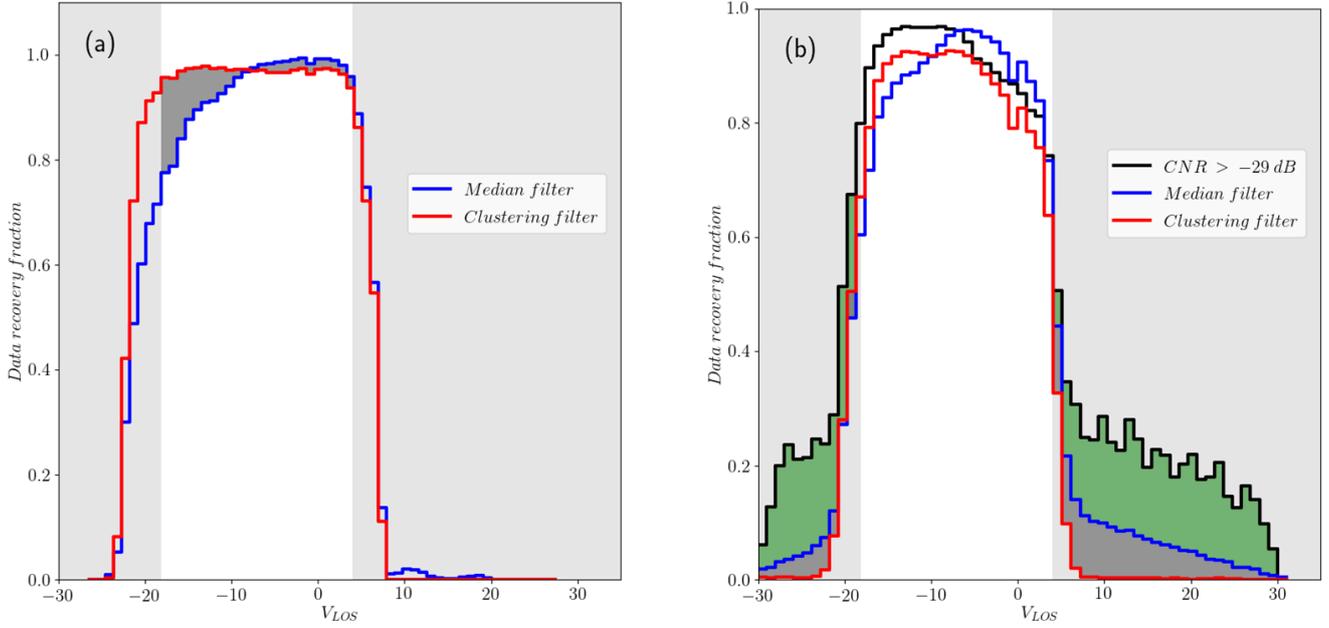


**Figure 9.** Histograms of the three performance indexes for the total number of synthetic scans (a) Both filters show similar spread but the clustering filter rejects a rather higher fraction of noise. (b) The higher recovery rate of the median-like filter, and its narrower distribution is superior than the clustering algorithm, the cost is acceptance of more contaminated observations (c) Both filters have similar mean values for  $\eta_{tot}$  around 0.9

## 285 6.2 Real data

The data set from the Balconies experiment presents advantages for the clustering filter, since the CNR value can be included as a feature in describing the data. Nevertheless, as mentioned already in section 2, we do not count on any reference to asses the performance of the filter apart from the radial wind speeds distribution of very reliable observations with CNR values within the range between -24 dB and -8 dB. As mentioned earlier, valid observations in this range might present a similar distribution.

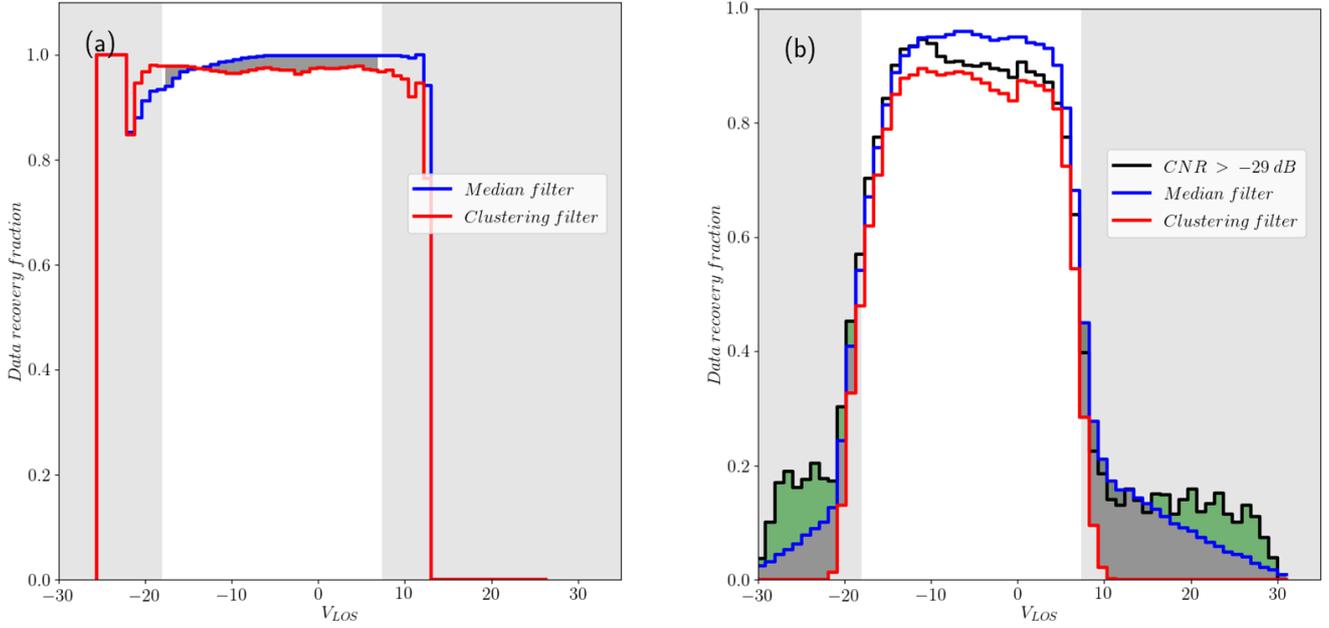
290 Figure 7 shows this distribution before filtering, shadowing the area of values of  $V_{LOS}$  that fall in the region beyond a 99.7% of the total probability or  $3\sigma$  limit, usually classified as outliers. Figures 10 and 11 show the recovery fraction for CNR, median-



**Figure 10.** Distribution of recovery fraction per wind speed bin for phase 1 of the experiment of (a) reliable observations ( $-24 < \text{CNR} < -8$ ) and (b) non reliable data ( $\text{CNR} < -24$  or  $\text{CNR} > -8$ ) for the three types of filter. The shadowed area in both graphs corresponds to the region where observations exceeded the 99.7% of probability (or  $3\text{-}\sigma$  limit) in the pdf of reliable observations. The darker shadowed areas highlights the additional fraction of extreme values non-filtered by the median-like and CNR filters, when the former uses the optimal input set  $n_r = 5$ ,  $n_\phi = 3$  and  $\Delta V_{LOS, threshold} = 2.33$  m/s.

like and clustering filters when applied on data in the reliable and non reliable CNR ranges for phases 1 and 2 of the Østerild experiment. Unlike the clustering filter, the CNR threshold and median-like filters show non negligible recovery rates beyond the  $3\text{-}\sigma$  limit, particularly significant in the former. This result is very much in line with the  $\eta_{noise}$  metric from the synthetic data. Within the  $3\text{-}\sigma$  range, the CNR and median-like filters perform slightly better than the clustering filter in terms of recovery fraction, in agreement with the results of  $\eta_{rec}$  in section 6.1. Even though this might compensate the fact that CNR threshold and median-like filters fail to filter out the major part of outliers, increasing the availability of measurements, this difference does not make the pdf of the filtered data more similar to the pdf of reliable data, as Table 4 shows via the metrics  $D_K$  and  $D_{KL}$ . According to this metric, the pdf of the data after the application of the clustering approach looks more statistically similar to reliable observations. This table also show  $D_K$  and  $D_{KL}$  of the non reliable data before filtering, which in all cases is improved, except for  $D_K$  for median and CNR threshold filters during phase 2.

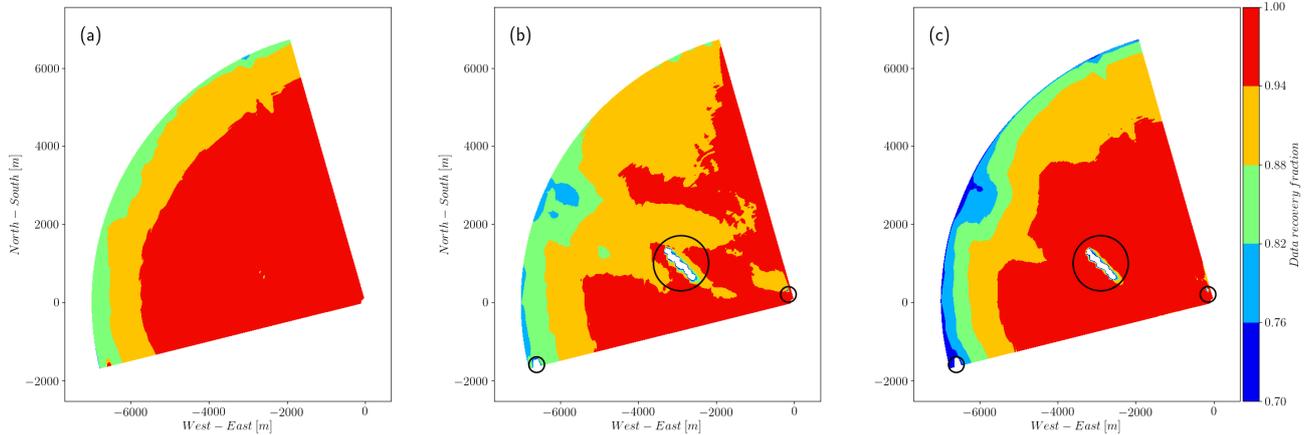
Figures 12, 13 and 14 show the performance of the three different filters in different regions of the scan, from respectively phase 1 and 2 of the experiment. When the spatial distribution of the recovery fraction is analyzed, we can see that the lowest values shown by the clustering filter are mostly located in the far region of the scan which, in general, presents low CNR values. The spatial recovery rate during phase 1 also show that the median-like and clustering filters are able to identify hard



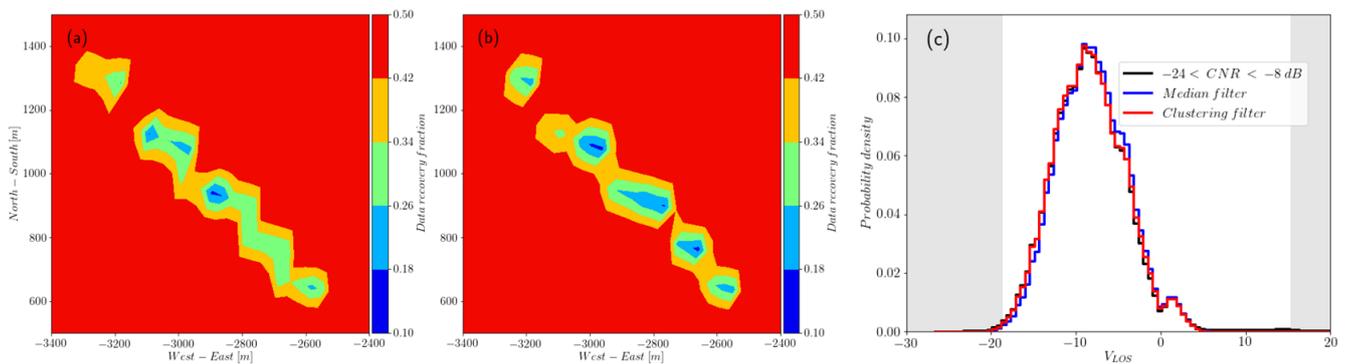
**Figure 11.** Distribution of recovery fraction per wind speed bin for phase 1 of the experiment of (a) reliable observations and (b) non reliable data for CNR, median and clustering filter. The shadowed area in both graphs corresponds to the region where observations exceed the  $3\text{-}\sigma$  limit in the pdf of reliable observations. Again, darker shadowed areas highlights the additional fraction of extreme values non-filtered by the median-like and CNR filters, when the former uses the optimal input set  $n_r = 5$ ,  $n_\phi = 3$  and  $\Delta V_{LOS, threshold} = 2.33$  m/s.

targets, which are also a source of bad observations. For scans recorded at 50 m above ground level in phase 1, back-scatter is affected by a group of seven turbines located approximately in the middle of the scanning area, with one turbine touching the end of the southern beams of the scan and a meteorological mast located very close to the lidar. Figure 13 shows a detail of the recovery rate associated with the flow in the vicinity of the turbines group, in which we can see that the clustering filter is able to identify better the turbine locations, recovering more data in the surroundings when compared to the median-like filter. The PDF of  $V_{LOS}$  in this area also show more similarities between the data filtered with the clustering algorithm and observations with CNR values in  $[-24, 8]$ .

Table 5 shows a summary of the additional data available when the  $\text{CNR} = -29$  dB threshold, the median-like and the clustering filters are applied instead of the more conservative and restrictive  $\text{CNR} = -24$  dB threshold filter. Additionally, this table shows the fraction of observations exceeding the  $3\text{-}\sigma$  limit that are recovered by the three filters. Even though the clustering filter shows a slightly lower fraction of additional data available when compared to the other filters, most of it comes from values within the  $3\text{-}\sigma$  region. Moreover the quality of the data recovered by the clustering approach seems to be higher when all these results are tested with the performance metrics defined in Section 5.



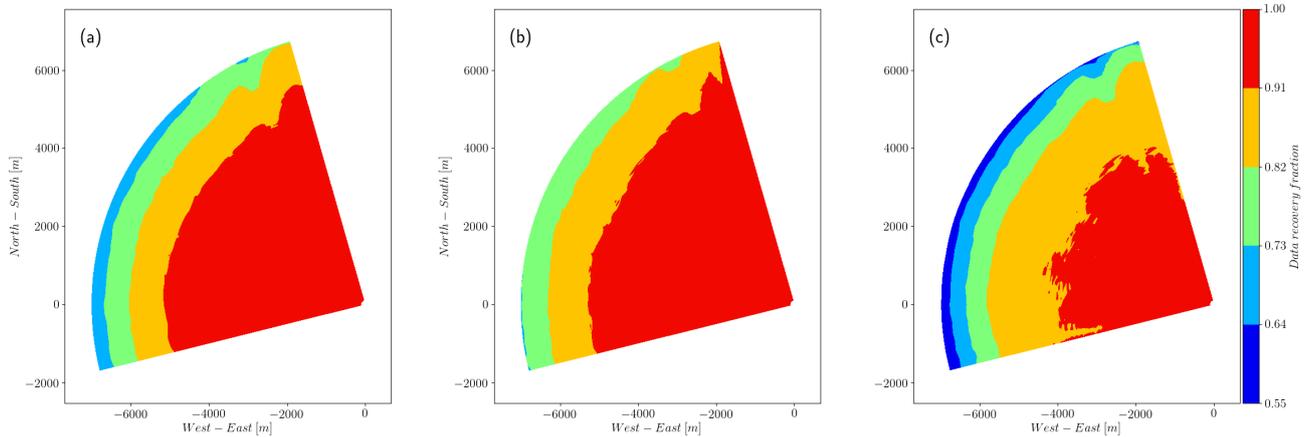
**Figure 12.** Total recovery fraction for phase 1 of the experiment. The noisy and far region of the scans show a high recovery, above 80%, for (a) the  $CNR > -29$  dB threshold filter and (b) the median-like filter and below 75% for (c) the clustering filter. Highlighted, it is possible to see three groups of hard targets (turbines and one meteorological mast, close to the lidar), which are identified by the median and clustering filter with recovery rates below 20%.



**Figure 13.** Detail of the recovery rate at the site of the turbines for (a) median filter and (b) clustering filter. The recovery is lower in the flow regime of the turbines cluster (there are 7 turbines in line) and higher in their surrounding for the clustering filter. Red denotes recovery rates of 0.5 or higher. (c) Probability density of  $V_{LOS}$  around the group of turbines

## 7 Discussion

320 The metrics introduced in section 5.1 attempt to evaluate two different capabilities of the filters: the quality and amount of the data recovered. In general these two metrics are in conflict, every time a high rate of noise will decrease the data recovery. The metric  $\eta_{tot}$  attempts to quantify their relative importance regarding the noise fraction, which in this study is distributed in a relatively wide range, but on average represents 20% of the total number of measurements per scan. The impact of the noise



**Figure 14.** The total recovery fraction of observations for phase 2 of the experiment. The noisy and far region of the scans show a high recovery, above 70% for (a) and (b), the CNR  $> -29$  dB threshold and median-like filters, respectively. The recovery decreases to 55% in the same region for the clustering filter, in line with the previous results, assuming that outliers (above the  $3\sigma$  limit) and noise are more likely to be located here.

**Table 4.** Results of pdf similarity test of reliable and non-reliable data after filtering. The CNR =  $-29$  dB threshold is also included (Gryning and Floors (2019))

Phase 1	$D_K$	$D_{KL}$
Non-reliable data before filtering	0.097	0.134
CNR threshold $> -29$ dB	0.045	0.109
Median filter	0.047	0.126
Clustering filter	0.037	0.105
Phase 2		
Non-reliable data before filtering	0.110	0.126
CNR threshold $> -29$ dB	0.114	0.052
Median filter	0.117	0.057
Clustering filter	0.103	0.045

fraction distribution on the performance of the filters was not explored, and variations on its dispersion and mean value might  
 325 be necessary. Regarding the synthetic scans, they do not allow the identification of outliers in the time domain because they  
 are time independent. A time evolving synthetic turbulence fields would be necessary to generate scans correlated in time and  
 enhance the self similarity of the data. This might improve the performance of the clustering approach and allow the addition  
 of a time dependence in the median-like filter, used already in Meyer Forsting and Troldborg (2016).

**Table 5.** Additional data recovered, relative to the amount of observations in the reliable range of CNR, and fraction of data recovered with values beyond the  $3\text{-}\sigma$  range.

Phase 1, $3\sigma$ quantiles = [-18.16, 3.96] m/s	Fraction of data recovered beyond $3\text{-}\sigma$	Additional data recovered
CNR threshold > -29 dB	27.1%	23.4%
Median filter	14.0%	23.1%
Clustering filter	8.6 %	22.1%
Phase 2, $3\sigma$ quantiles = [-18.08, 7.35] m/s		
CNR threshold > -29 dB	16.5%	40.4%
Median filter	12.6%	42.4%
Clustering filter	3.2%	38.1%

The synthetic wind fields used here do not consider the presence of hard targets. These anomalies in the wind field are  
 330 observed by lidars as points with high CNR values and abnormal  $V_{LOS}$ . Assessing the performance of the filters in detecting  
 such anomalies needs a more realistic model of the pulsed lidar. This lidar simulator would allow the generation of information  
 normally available in real lidar measurements, like CNR, and the spread in the power spectra of the heterodyne signal,  $S_b$ . This  
 additional information will benefit the performance assessment of the clustering filter and the simulation of hard targets. A  
 more realistic lidar model was already implemented by Brousmiche et al. (2007), which can be used to explore further these  
 335 aspects of the filtering process.

The data set analyzed from the Balconies experiment corresponds to horizontal scans at 50 and 200 m above the ground  
 level, limiting the analysis to one scanning pattern. Different scanning patterns, in vertical and horizontal planes, as well as wind  
 fields over different topography would make this analysis more general, thus shedding light on the capabilities of the filters here  
 presented. This is specially critical regarding the median-like filter, which might require again a sensitivity analysis to select  
 340 proper parameters that adapt to different scanning patterns and turbulence field characteristics. So far,  $\Delta V_{LOS, threshold}$  showed  
 a dependence on the  $L$  and  $\alpha\varepsilon^{2/3}$  parameters during the sensitivity analysis presented in Section 6.1. Larger fluctuations in the  
 $V_{LOS}$  field, whether they come from larger turbulent structures or higher turbulence energy or both, will need a larger value  
 of  $\Delta V_{LOS, threshold}$  to avoid the rejection of good measurements. Range Height Indicator (RHI) scanning patterns can pose  
 the challenge of strong vertical shear and small turbulent structures that will need to reduce the window size  $n_r$  and  $n_\phi$  for the  
 345 median-like filter, and the selection of a different set of features (or a new definition for  $\Delta V_{LOS}$ ) for the clustering filter, in  
 order to keep reliable observations from being filtered out.

Regarding feature selection, the clustering filter could consider the spectral spreading of the heterodyne signal,  $S_b$  and time  
 variation of  $V_{LOS}$ , in addition to features already used in this work to characterize and distinguish better cluster of good  
 measurements. Nevertheless, due to the Euclidean distance definition, additional dimensions will make the data more sparse  
 350 in higher dimensions, making it necessary to use more data points per filtering step (here we used only 3 scans at a time) to

avoid the identification of good observations as spread, low density noise. It is because of this that the application of a feature selection method might be necessary (Chandrashekar and Sahin, 2014).

Using the statistical distances  $D_K$  and  $D_{KL}$  as a metric for the filter performance might not be totally correct. At range gates far from the lidar, the distance between beams increases, as well as the area covered by the accumulation of spectral information in azimuth direction. Averaging  $V_{LOS}$  over larger areas as we move forward through each beam, might affect the statistics and the PDF of  $V_{LOS}$  (specially its spread) in the outer region of the scan. The fact that this region is where we usually find the non reliable measurements group, may make the PDFs of reliable and non reliable observations somewhat different. These possible deviations need to be investigated further.

The selection of features and the amount of scans put together per filtering step/iteration could also be automatized, using feature selection methods. Nevertheless, this would make the clustering filter more complex in its implementation and more computationally expensive, which is the main disadvantage of this approach compared to the median-like filter. Very efficient median filters can achieve a computational complexity up to  $\mathcal{O}(n)$ , with  $n$  being the number of observations in the data set. Depending on the data structure, DBSCAN shows a computational complexity from  $\mathcal{O}(n \log(n))$  to  $\mathcal{O}(n^2)$ . If the distance between points is in general smaller than  $\varepsilon$ , the first limit can be achieved, but clusters with different densities makes the algorithm less efficient. In the data analyzed here, having clusters with different densities is not an issue. Nevertheless, for non homogeneous flows, scans might persistently show regions with  $V_{LOS}$ , CNR or other feature with noticeable different values, may need to revisit the clustering algorithm used and implement a  $\varepsilon$ -independent clustering approach, like OPTICS (Ankerst et al., 1999) for instance.

## 8 Conclusions

The CNR threshold filtering has been the common approach to retrieve reliable observations from lidars measurements. In this work we compared this approach against two alternative techniques: a median-like filter, based on the assumption of smoothness of the wind field, hence, in the smoothness of the radial wind speed observed by a wind lidar, and a clustering filter, based in the assumption of self-similarity of the observations captured by the wind lidar and the possibility of clustering them in groups of good data and noise. A controlled test was carried out on the last two approaches, using a simple lidar simulator that sampled scans from synthetic wind fields, later contaminated with procedural noise. The results indicate that the clustering filter is capable of detecting more added noise than the median-like filter, at a good recovery rate of non contaminated data. When the three filters are tested on real data, the clustering approach shows a better performance on identifying abnormal observations, increasing the data availability between 22% and 38% and reducing the recovery of abnormal measurements between 70 and 80% when compared to a CNR threshold. This is an important result, because increases the spatial coverage of the data which can be used later for wind field reconstruction and wind data analysis, specially in the far region of the scan, that covers the largest measured area.

Even though the median-like filter is computationally efficient, it needs an optimal definition of input parameters, which are dependent on the turbulence characteristics of the wind field. The clustering filter is more robust in this sense, because it is

capable of automatically adapt its parameters to the structure of the data. This is a step forward to a more robust and automated  
385 processing of data from lidars, which ideally should be independent of the turbulence characteristics of the measured wind  
field or the scanning pattern used.

*Code and data availability.* The synthetic data and code is available at [https://github.com/lalcayag/Lidar\\_filtering](https://github.com/lalcayag/Lidar_filtering). Real data can be found  
in Simon and Vasiljevic (2018)

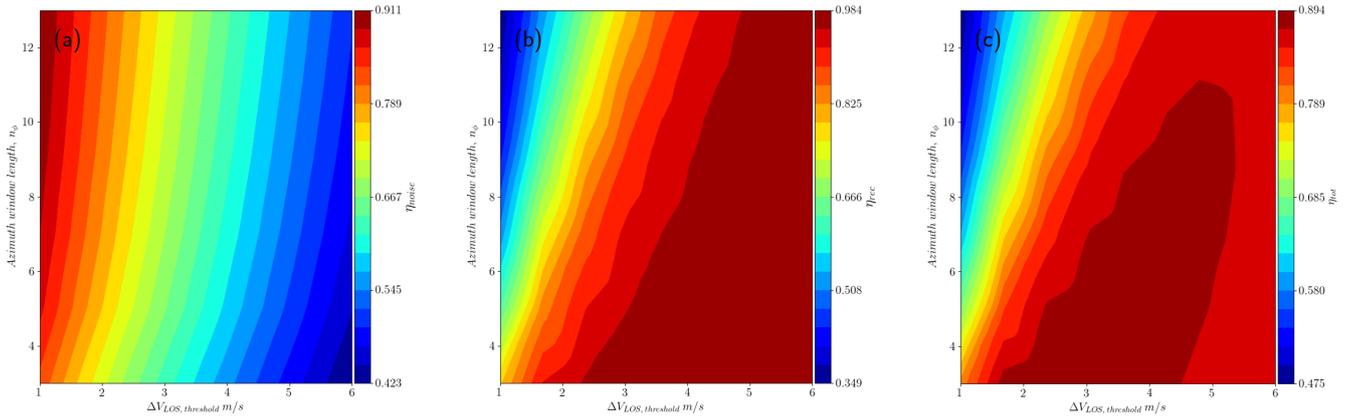
*Author contributions.* The author Leonardo Alcaayaga implemented the analysis on synthetic and real data and wrote all sections of this  
390 manuscript

*Competing interests.* The author declares not having any competing interest

*Acknowledgements.* I would like to thank Ioanna Karagali for comments and guide through the data from the Østerild Balconies experiments,  
Robert Menke for the first version of the median-like filter and Ebba Dellwik, Nikola Vasiljevic, Gunner Larsen, Mark Kelly and Jakob Mann  
for the very useful and clarifying feedback during the analysis and writing process.

## 395 **Appendix A: Sensitivity analysis on median-like filter parameters**

Figure A1 shows contours that present the most optimal value for  $\eta_{tot}$  among all possible values of  $\Delta V_{LOS, threshold}$  and  $n_\phi$ ,  
for  $n_r = 5$ , the optimal window size in the radial direction. Large  $\Delta V_{LOS, threshold}$  results in large  $\eta_{recov}$  but poor results  
for  $\eta_{noise}$  and the opposite for values of the threshold, as expected. The metric  $\eta_{tot}$  then becomes relevant to determine the  
optimal combination of parameters. From the contours it is possible to see that the performance in terms of the  $\eta_{tot}$  metric  
400 is less sensitive to  $n_\phi$  than  $\Delta V_{LOS, threshold}$ . Even though the results here show average metrics for all the scans simulated,  
the optimal value of  $\Delta V_{LOS, threshold}$  increases with the turbulence energy and length scale parameters, which is problematic,  
because it requires previous knowledge of turbulence characteristics that usually are not available before reconstruction, and  
more important, data filtering.



**Figure A1.** Contours of performance metrics for  $n_r = 5$  over the  $\Delta V_{LOS, threshold} - n_\phi$  space. Each point in the contour plot corresponds to the mean value of (a)  $\eta_{noise}$ , (b)  $\eta_{rec}$  and (c)  $\eta_{tot}$  among all the 4305 synthetic scans filtered. The optimal value corresponds to  $n_r = 5$ ,  $n_\phi = 3$  and  $\Delta V_{LOS, threshold} = 2.33$  m/s

## References

- 405 Ankerst, M., Breunig, M. M., Peter Kriegel, H., and Sander, J.: OPTICS: Ordering Points To Identify the Clustering Structure, in: Proc. ACM SIGMOD'99 Int. Conf. on Management of Data, pp. 49–60, ACM Press, 1999.
- Backer, E.: Computer-assisted Reasoning in Cluster Analysis, Prentice Hall International (UK) Ltd., Hertfordshire, UK, UK, 1995.
- Banakh, V. A. and Smalikho, I. N.: “Estimation of the Turbulence Energy Dissipation Rate from the Pulsed Doppler Lidar Data, Atmos. Ocean. Opt., 10 (12), 957–965, 1997.
- 410 Beck, H. and Kühn, M.: Dynamic Data Filtering of Long-Range Doppler LiDAR Wind Speed Measurement., Remot Sens, 9(6), 561, <https://doi.org/https://doi.org/10.3390/rs9060561>, 2017.
- Brousmiche, S., Bricteux, L., Sobieski, P., Macq, B., and Winckelmans, G.: Numerical simulation of a heterodyne Doppler LIDAR for wind measurement in a turbulent atmospheric boundary layer, in: 2007 IEEE International Geoscience and Remote Sensing Symposium, <https://doi.org/10.1109/IGARSS.2007.4423420>, 2007.
- 415 Burger, W. and Burge, M. J.: Digital Image Processing - An Algorithmic Introduction using Java., Texts in Computer Science, Springer, 2008.
- Cariou, J.: Remote Sensing for Wind Energy, chap. Pulsed lidars, pp. 131–148, DTU Wind Energy, Denmark, 2015.
- Chandrashekar, G. and Sahin, F.: A Survey on Feature Selection Methods, Computers and Electrical Engineering, 40, 16–28, <https://doi.org/10.1016/j.compeleceng.2013.11.024>, 2014.
- 420 Ester, M., Kriegel, H.-P., Sander, J., and Xu, X.: A Density-based Algorithm for Discovering Clusters a Density-based Algorithm for Discovering Clusters in Large Spatial Databases with Noise, in: Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, KDD'96, pp. 226–231, Portland, Oregon, 1996.
- Gryning, S.-E. and Floors, R.: Carrier-to-Noise-Threshold Filtering on Off-Shore Wind Lidar Measurements, Sensors, 19, <https://doi.org/10.3390/s19030592>, 2019.

- 425 Gryning, S.-E., Floors, R., Peña, A., Batchvarova, E., and Brümmner, B.: Weibull Wind-Speed Distribution Parameters Derived from a Combination of Wind-Lidar and Tall-Mast Measurements Over Land, Coastal and Marine Sites., *Bound-Lay Meteorol*, 159(2), 329, <https://doi.org/https://doi.org/10.1007/s10546-015-0113-x>, 2016.
- Huang, T., Yang, G., and Tang, G.: A fast two-dimensional median filtering algorithm, *IEEE T. Acoust. Speech*, 27, 13–18, <https://doi.org/10.1109/TASSP.1979.1163188>, 1979.
- 430 Karagali, I., Mann, J., Dellwik, E., and Vasiljević, N.: New European Wind Atlas: The Østerild balconies experiment, *J Phys Conf Ser*, 1037, 052029, <https://doi.org/10.1088/1742-6596/1037/5/052029>, 2018.
- Kolmogorov, A.: Sulla determinazione empirica di una legge di distribuzione, *Inst. Ital. Attuari, Giorn.*, 4, 83–91, <https://ci.nii.ac.jp/naid/10010480527/en/>, 1933.
- Kullback, S. and Leibler, R. A.: On Information and Sufficiency, *Ann. Math. Statist.*, 22, 79–86, 1951.
- 435 MacQueen, J.: Some methods for classification and analysis of multivariate observations, in: *Proceedings Fifth Berkeley Symp. on Math. Statist. and Prob.*, vol. 1: Statistics, pp. 281–297, Berkeley, California, 1967.
- Mandelbrot, B. B.: *The fractal geometry of nature*, W. H. Freeman and Comp., New York, 1983.
- Mann, J.: The spatial structure of neutral atmospheric surface-layer turbulence, *J Fluid Mech.*, 273, 141–168, <https://doi.org/10.1017/S0022112094001886>, 1994.
- 440 Mann, J.: Wind field simulation, *Probabilist. Eng. Mech.*, 13, 269–282, [https://doi.org/https://doi.org/10.1016/S0266-8920\(97\)00036-2](https://doi.org/https://doi.org/10.1016/S0266-8920(97)00036-2), 1998.
- Mann, J., Angelou, N., Arnqvist, J., Callies, D., Cantero, E., Arroyo, R. C., Courtney, M., Cuxart, J., Dellwik, E., Gottschall, J., et al.: Complex terrain experiments in the new european wind atlas, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 375, 20160101, 2017.
- 445 Menke, R., Vasiljević, N., Wagner, J., Oncley, S. P., and Mann, J.: Multi-lidar wind resource mapping in complex terrain, *Wind Energy Science Discussions*, 2019, 1–21, <https://doi.org/10.5194/wes-2019-85>, <https://wes.copernicus.org/preprints/wes-2019-85/>, 2019.
- Meyer Forsting, A. and Troldborg, N.: A finite difference approach to despiking in-stationary velocity data - tested on a triple-lidar, *J Phys Conf Ser*, 753, 072017, <https://doi.org/10.1088/1742-6596/753/7/072017>, 2016.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., 450 Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E.: Scikit-learn: Machine Learning in Python, *J. Mach. Learn. Res.*, 12, 2825–2830, 2011.
- Perlin, K.: Noise hardware, In *Real-Time Shading SIGGRAPH Course Notes*, 2001.
- Simon, E. and Vasiljevic, N.: Østerild Balconies Experiment (Phase 2), <https://doi.org/10.11583/DTU.7306802.v1>, [https://data.dtu.dk/articles/\\_sterild\\_Balconies\\_Experiment\\_Phase\\_2\\_/7306802](https://data.dtu.dk/articles/_sterild_Balconies_Experiment_Phase_2_/7306802), 2018.
- 455 Smalikho, I. N. and Banakh, V. A.: Accuracy of estimation of the turbulent energy dissipation rate from wind measurements with a conically scanning pulsed coherent Doppler lidar. Part I. Algorithm of data processing, *Atmos. Ocean. Opt.*, 26, 404–410, <https://doi.org/10.1134/S102485601305014X>, 2013.
- Stawiarski, C., Träumner, K., Knigge, C., and Calhoun, R.: Scopes and Challenges of Dual-Doppler Lidar Wind Measurements—An Error Analysis, *Journal of Atmospheric and Oceanic Technology*, 30, 2044–2062, <https://doi.org/10.1175/JTECH-D-12-00244.1>, <https://doi.org/10.1175/JTECH-D-12-00244.1>, 2013.
- 460 Vasiljevic, N., Lea, G., Courtney, M., Cariou, J.-P., Mann, J., and Mikkelsen, T.: Long-Range WindScanner System, *Remote Sensing*, 8, <https://doi.org/10.3390/rs8110896>, 2016.

Vasiljević, N., L. M. Palma, J. M., Angelou, N., Carlos Matos, J., Menke, R., Lea, G., Mann, J., Courtney, M., Frölen Ribeiro, L., and M. G. C. Gomes, V. M.: Perdigão 2015: methodology for atmospheric multi-Doppler  
465 lidar experiments, *Atmospheric Measurement Techniques*, 10, 3463–3483, <https://doi.org/10.5194/amt-10-3463-2017>, 2017.